

Microsoft
WinHEC
2005

Multi-Core Processor Technology – Maximizing CPU Performance in a Power-Constrained World

Paul Teich
Business Strategy
CPG Server/Workstation
paul.teich@amd.com
AMD



The Issues

- Silicon designers can choose a variety of methods to increase processor performance
- Commercial end-customers are demanding:
 - More capable systems with more capable processors
 - That new systems stay within their existing power/thermal infrastructure
- Processor frequency and power consumption seem to be scaling in lockstep
- How can the industry-standard PC and Server industries stay on our historic performance curve without burning a hole in our motherboards?
- This session is not about process technology...

Session Outline

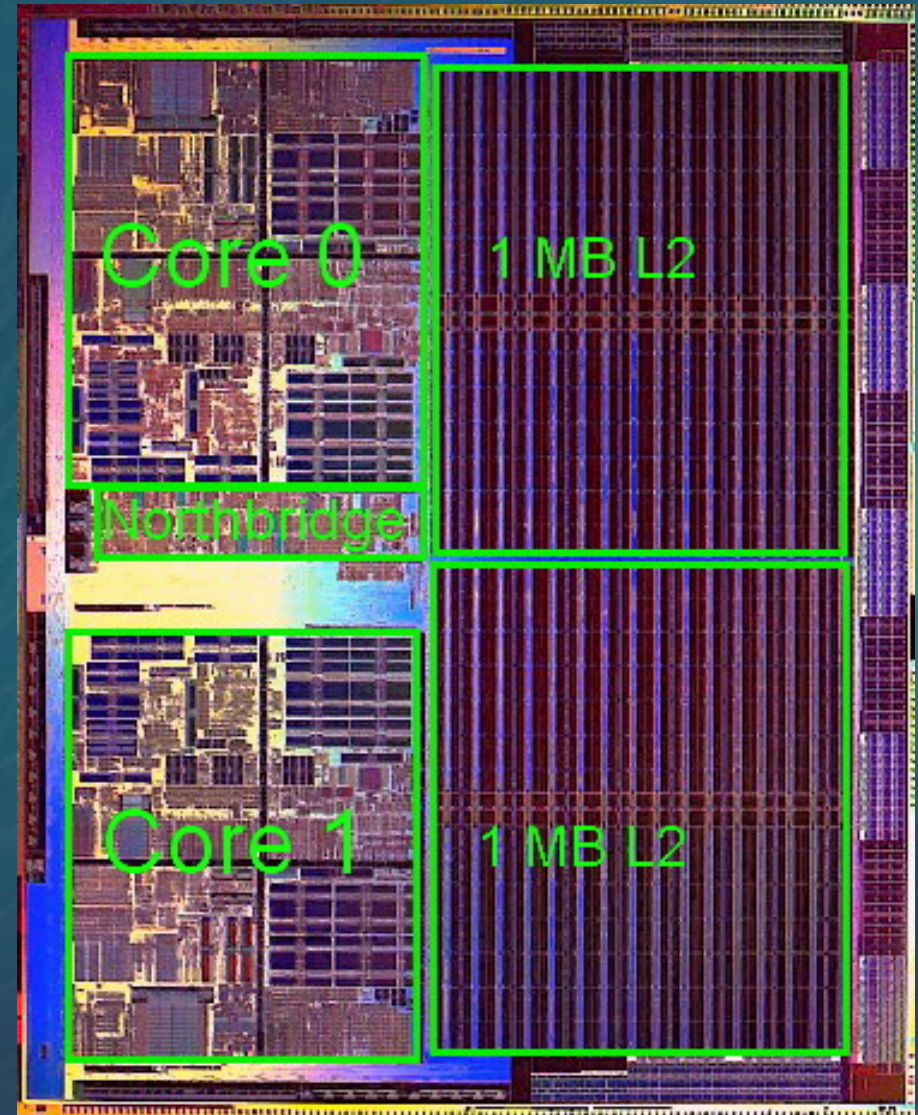
- Definition – What is a processor?
- Core Design
- System Architecture
- Manufacturing, Power, and Thermals
- Multi-Core Processor Architecture
- Performance Impacts

What is a Processor?

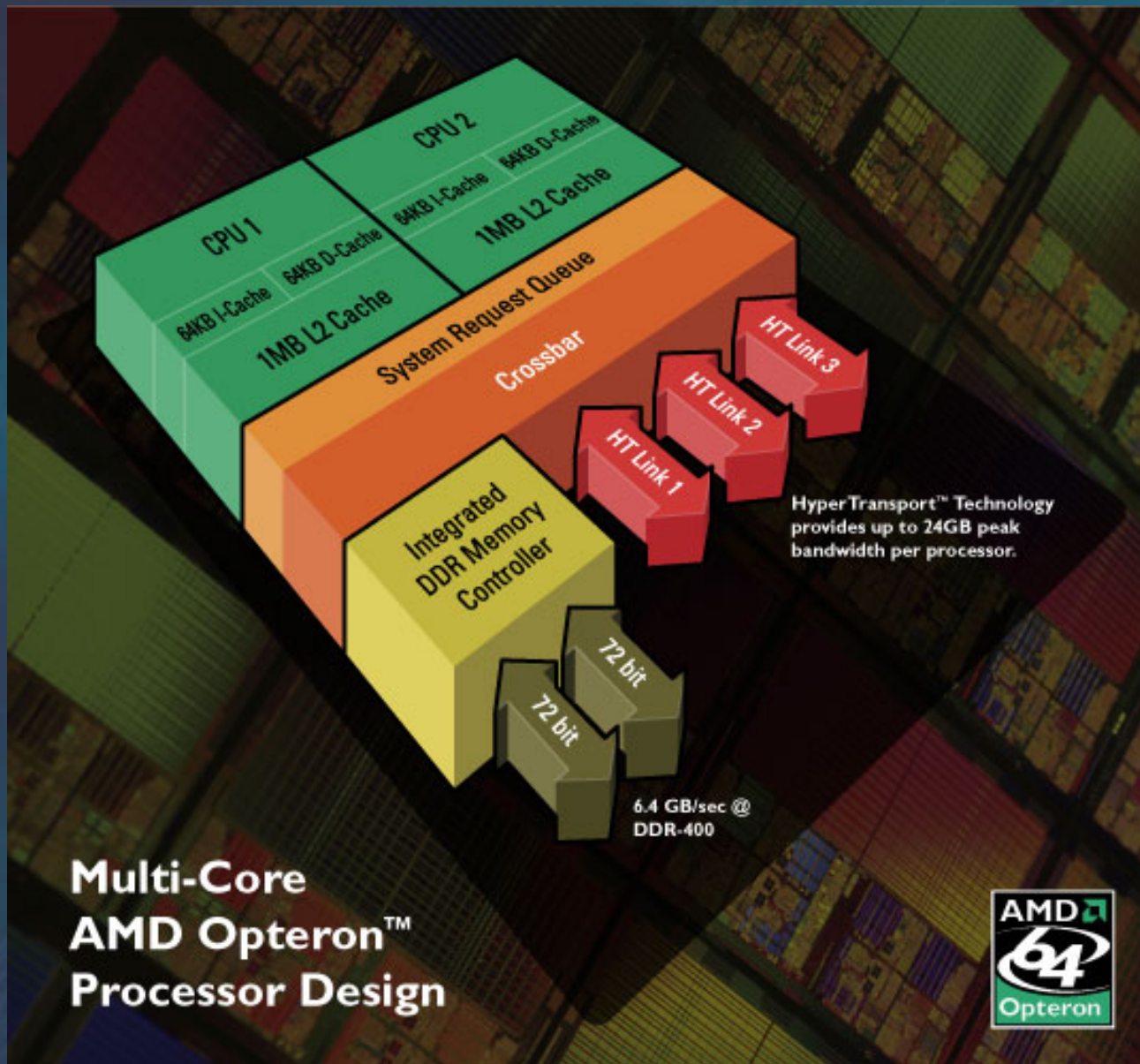
- A single chip package that fits in a socket
- ≥ 1 cores (not much point in < 1 core...)
 - Cores can have functional units, cache, etc. associated with them, just as today
 - Cores can be fast or slow, just as today
- Shared resources
 - More cache
 - Other integration – northbridge, memory controllers, high-speed serial links, etc.
- One system interface no matter how many cores
 - Number of signal pins doesn't scale with number of cores

A Representative Multi-Core Processor

- Dual-core AMD Opteron™ processor is 199mm² in 90nm
- Single-core AMD Opteron processor is 193mm² in 130nm

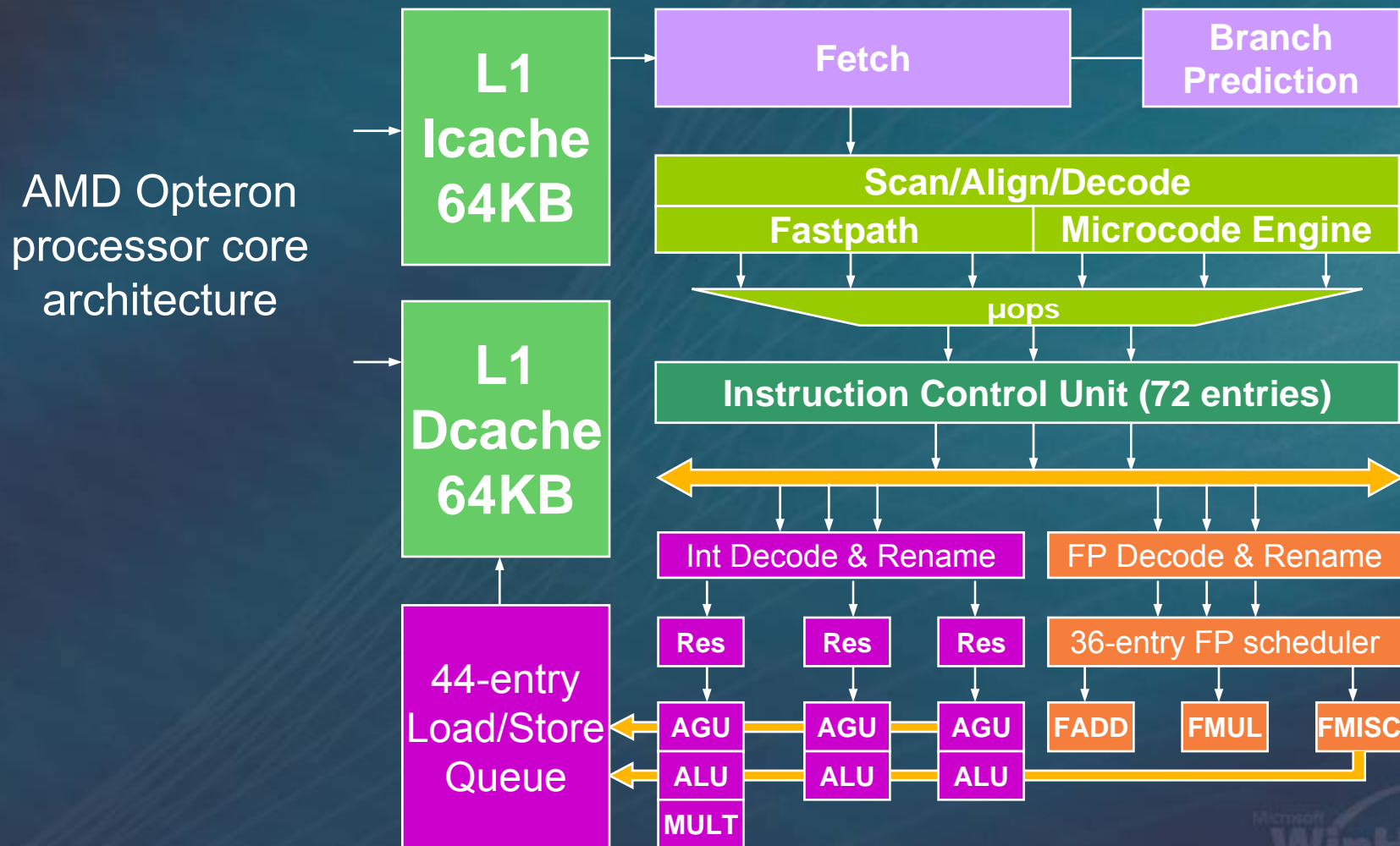


Multi-Core Processor Architecture



Core Design

- Frequency
 - Is only as good as the rest of the core architecture



Core Design

- Functional units
 - Superscalar is known territory
 - Diminishing returns for adding more functional blocks
 - Alternatives like VLIW have been considered and rejected by the market
 - Single-threaded architectural performance is pegged
- Data paths
 - Increasing bandwidth between functional units in a core makes a difference
 - Such as comprehensive 64-bit design, but then where to?

Core Design

- Pipeline

- Deeper pipeline buys frequency at expense of increased cache miss penalty and lower instructions per clock
- Shallow pipeline gives better instructions per clock at the expense of frequency scaling
- Max frequency per core requires deeper pipelines
- Industry converging on middle ground...9 to 11 stages
 - Successful RISC CPUs are in the same range

- Cache

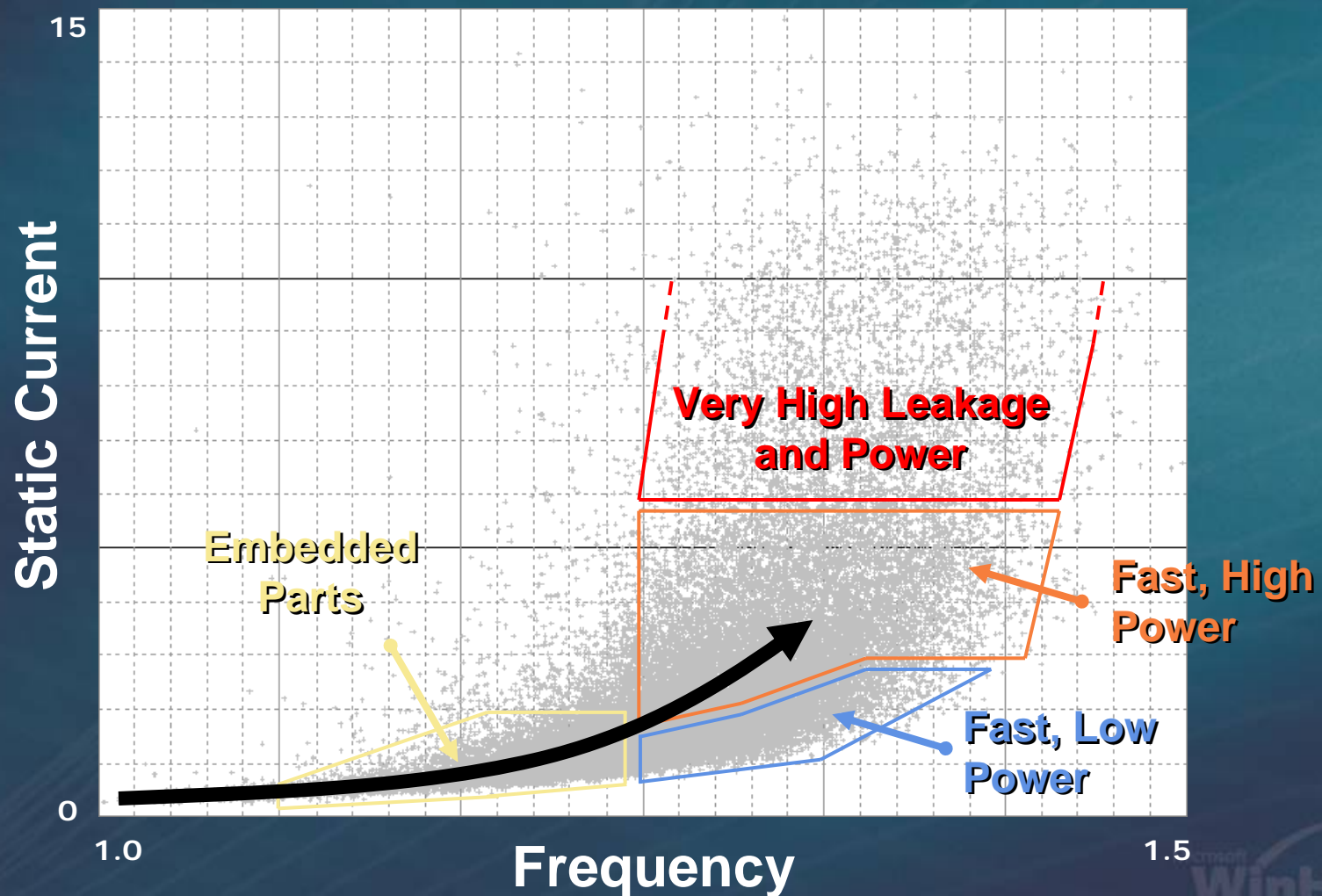
- Cache size buys performance at expense of die size, it's a direct hit to manufacturing cost
- Deep pipeline cache miss penalties are reduced by larger caches
- Not always the best match for shallow pipeline cores, as cache misses penalties are not as steep

Manufacturing

- Moore's Law isn't dead, more transistors for everyone!
 - But...it doesn't really mention scaling transistor power
- Chemistry and physics at nano-scale
 - Stretching materials science
 - Voltage doesn't scale yet
 - Transistor leakage current is increasing
- As manufacturing economies and frequency increase, power consumption is increasing disproportionately
- There are no process or architectural quick-fixes

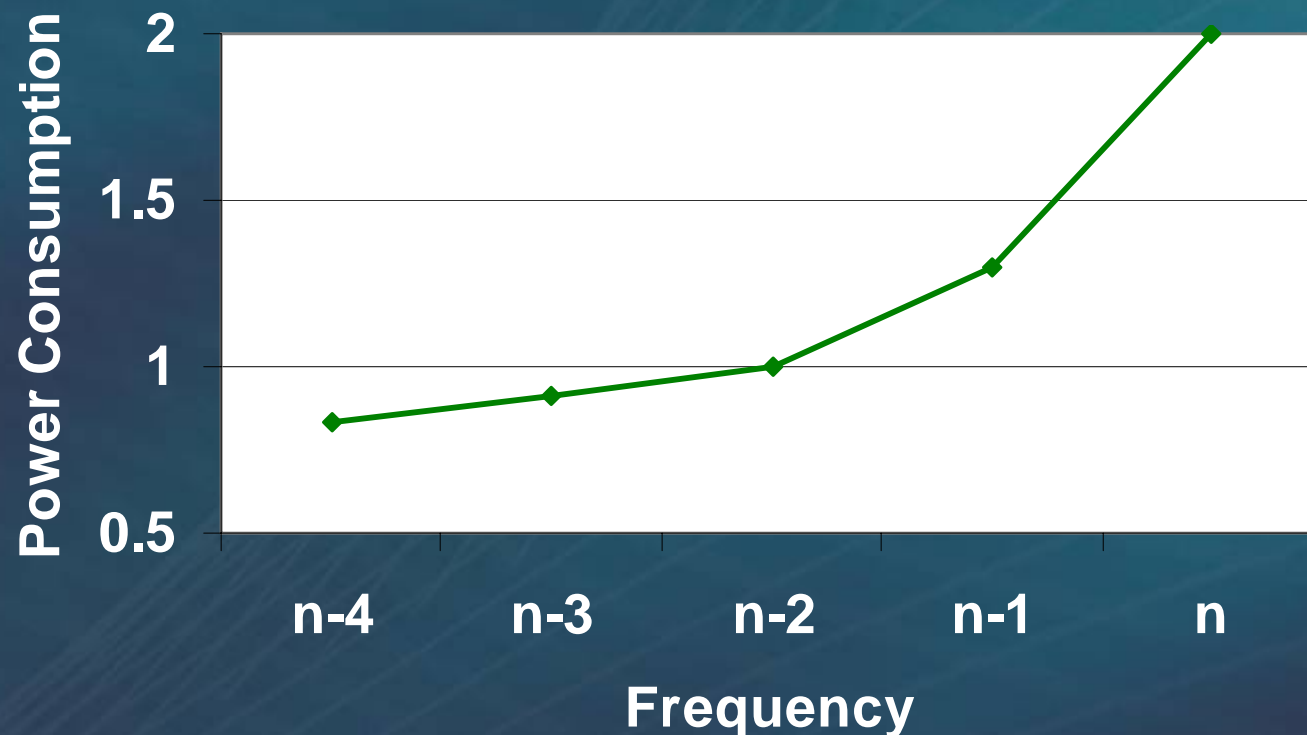
Static Current vs. Frequency

Non-linear as processors approach max frequency



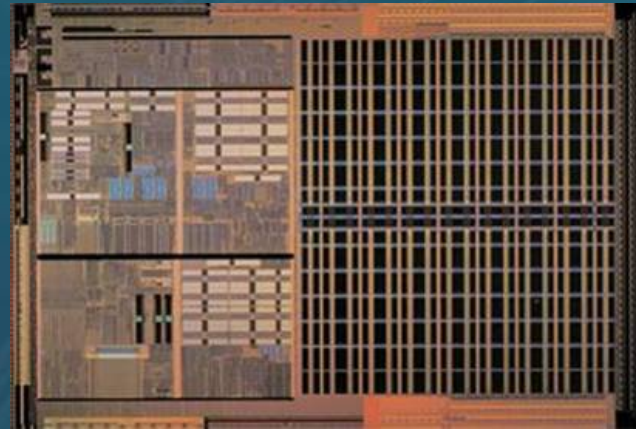
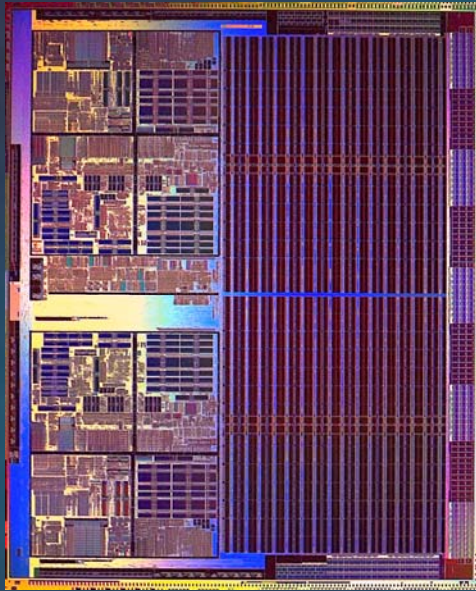
Power vs. Frequency

In AMD's process for 200MHz frequency steps, two steps back on frequency cuts power consumption by half...or lets you put twice the transistors on the die for the same power consumption as top frequency...



(Gross relative numbers summarized from a mountain of real data)

Thermal Density Decreases

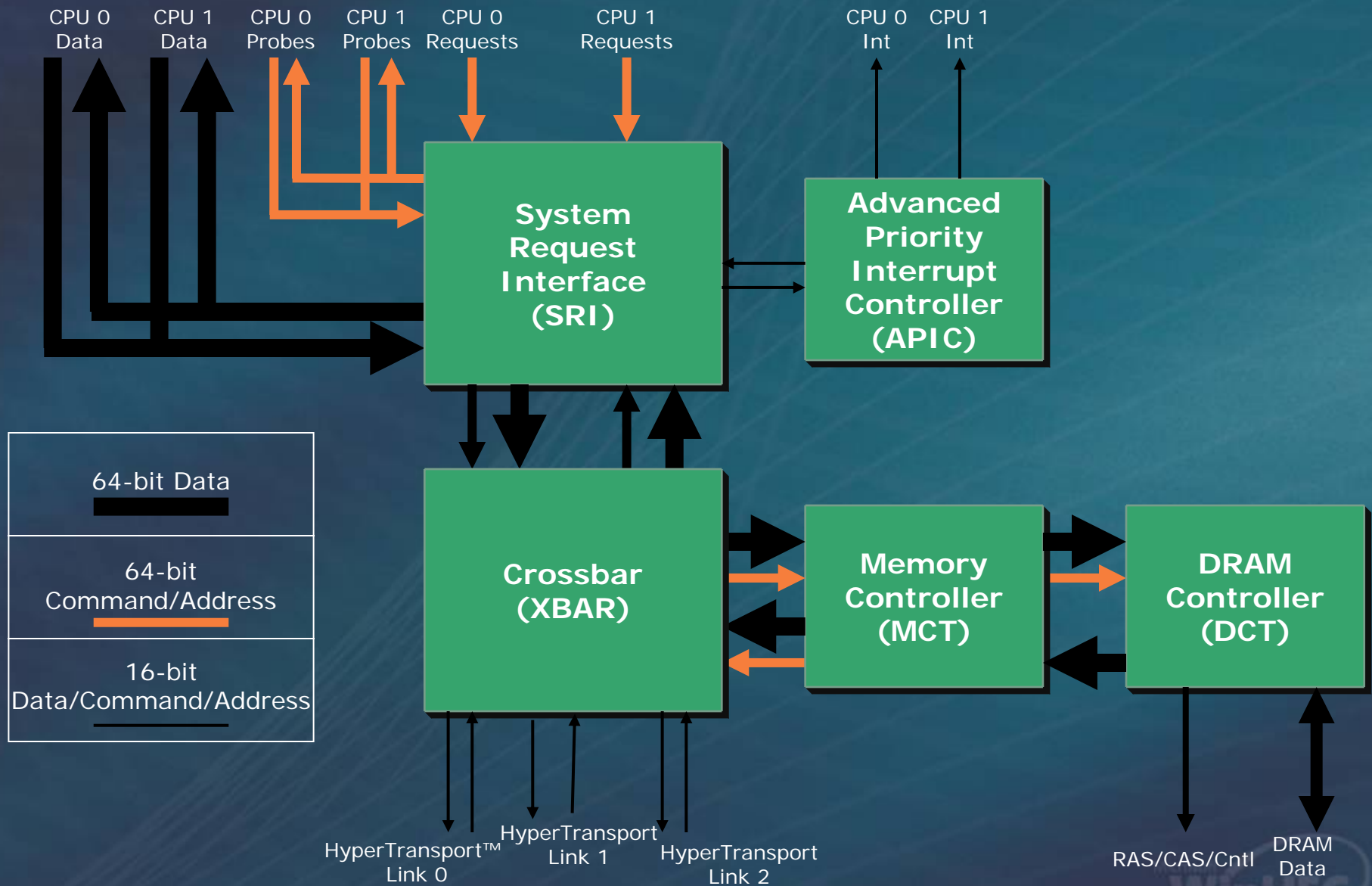


- Hot spots
 - Twice as many as in single-core
 - Farther apart than in single-core
 - With freq delta, cooler than in single-core
- Θ_{CA} same for single-core at n and dual-core at $n-2$
 - Larger die spreads heat more evenly in package
 - Use identical heat sink, slightly better cooling with dual-core
 - Works for this processor generation and next, Θ_{CA} changes over major generations
- Thermal diode accuracy becomes an issue with dual-core...

Multi-Core Processor Architecture

- Why integrate?
 - Most functions are really small compared to the cores and cache
 - All integrated logic runs at core frequency regardless of I/O speeds
- What to integrate?
 - Northbridge crossbar switch is key to integrating anything else
 - Memory controller to reduce memory latency and further reduce the need for cache
 - High-speed serial links for system I/O
- What not to integrate?
 - Most southbridge functions
 - Graphics

AMD Opteron Processor Integrated Northbridge

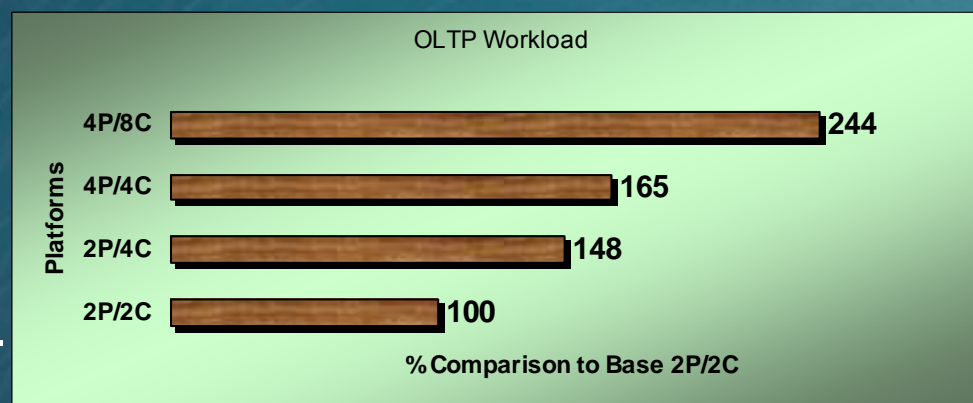
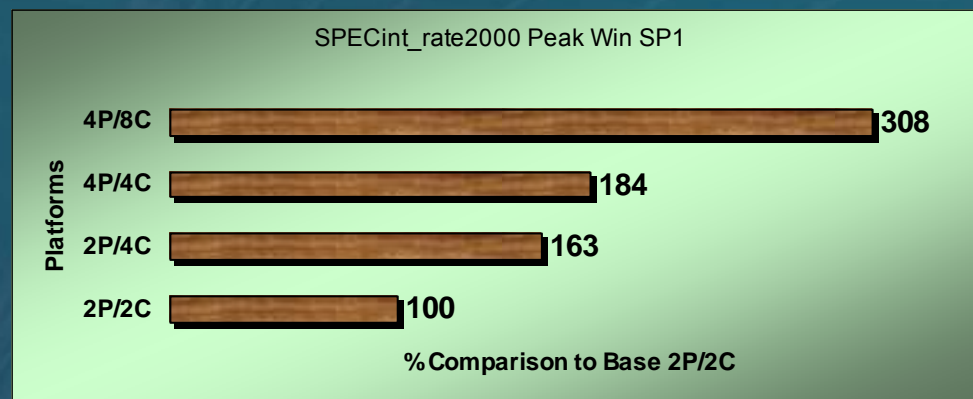


Multi-Core: Where Processor and System Collide

- Scales performance
 - Dedicated resources for two simultaneous threads
 - Multiple cores will contend for memory and I/O bandwidth
 - Northbridge is the bottleneck
 - Integrating northbridge eliminates most of bottleneck
 - Cores, cache and northbridge must be balanced for optimal performance
- More aggregate performance for:
 - Multi-threaded apps
 - Transactions – many instances of same app
 - Multi-tasking
- Thread scheduling handled by OS
 - BIOS notifies Windows of thread execution resources

Early Benchmark Estimates

- Decoder
 - 2P/2C – 2 proc. single-core
 - 4P/4C – 4 proc. single-core
 - 2P/4C – 2 proc. dual-core
 - 4P/8C – 4 proc. dual-core
- Frequencies
 - Single-core = 2.4GHz
 - Dual-core = 2.0GHz
- Identical system configs
 - Memory, disks, network, etc.
 - Early dual-core validation system used, different motherboards



SPEC and the benchmark name SPECint are registered trademarks of the Standard Performance Evaluation Corporation. SPEC scores for AMD Opteron Model 270 and 870 based systems are estimated.

Call To Action

- Most application software doesn't need to do anything to benefit from dual-core
- Be aware that, for a processor within a given power envelope:
 - Fewer cores will clock faster than more cores
 - Single-threaded performance-sensitive applications
 - More cores will out-perform fewer cores for
 - Multi-threaded applications
 - Multi-tasking response times
 - Transaction processing
- Processor architecture impacts multi-core performance
 - Process technology is only the ante
 - Integration enables a balanced high-performance architecture

Additional Resources

WinHEC Presentations:

- x86 Everywhere, Chris Herring, AMD

● Web Resources:

- AMD <http://www.amd.com/>
- AMD Multi-Core: <http://www.amd.com/multicore>
- AMD Opteron™ Processor Tech Docs: http://www.amd.com/us-en/Processors/TechnicalResources/0,,30_182_739_9003_00.html
- AMD Multi-Core White Paper: http://www.amd.com/us-en/assets/content_type/DownloadableAssets/MC_Whitepaper_vFINAL.pdf
- HyperTransport™ Consortium: <http://www.hypertransport.org/>

questions

Microsoft®

Your potential. Our passion.™

© 2005 Microsoft Corporation. All rights reserved.

This presentation is for informational purposes only. Microsoft makes no warranties, express or implied, in this summary.